



ELSEVIER

Rational approximation preconditioners for sparse linear systems

Philippe Guillaume^a, Yousef Saad^{b,1}, Masha Sosonkina^{c,*}

^aUMR MIP 5640, Département de Mathématiques, INSA, Complexe Scientifique de Rangueil, 31077 Toulouse Cedex, France

^bDepartment of Computer Science and Engineering, University of Minnesota, 200 Union Street S.E., Minneapolis, MN 55455, USA

^cDepartment of Computer Science, University of Minnesota-Duluth, 320 Heller Hall, 1114 Kirby Drive, Duluth, MN 55812-2496, USA

Received 12 March 2002; received in revised form 3 February 2003; accepted 18 March 2003

Abstract

This paper presents a class of preconditioning techniques which exploit rational function approximations to the inverse of the original matrix. The matrix is first shifted and then an incomplete LU factorization of the resulting matrix is computed. The resulting factors are then used to compute a better preconditioner for the original matrix. Since the incomplete factorization is made on a shifted matrix, a good LU factorization is obtained without allowing much fill-in. The result needs to be extrapolated to the nonshifted matrix. Thus, the main motivation for this process is to save memory. The method is useful for matrices whose incomplete LU factorizations are poor, e.g., unstable.

© 2003 Published by Elsevier Science B.V.

Keywords: ■; ■; ■

1. Introduction

Rational approximation preconditioners are targeted at extremely ill-conditioned linear systems. Examples of such systems are those that arise in the modeling of thin Shells. These systems tend to

* Corresponding author.

E-mail addresses: guillaum@gmm.insa-tlse.fr (P. Guillaume), saad@cs.umn.edu (Y. Saad), masha@d.umn.edu (M. Sosonkina).

¹ This work was supported in part by the U.S. Army Research Office under grant DAAD19-00-1-0485, and in part by the Minnesota Supercomputer Institute.

² This work was supported by Michelin Americas Research and Development Corporation.

1 be very difficult to solve by iterative methods despite the fact that they are symmetric positive definite.
 2 One source of difficulty is that the quality of incomplete LU (in this case Cholesky) factorizations for
 3 such matrices can be so poor that they become ineffective. It is tempting to simply shift the matrix
 4 A by a scalar α and extract the preconditioning for $A + \alpha I$ which is then used for preconditioning
 5 the original matrix, see, e.g., [8]. This by itself may not be sufficient.

6 A modification of this idea leads to a more effective technique. This modification consists of
 7 exploiting a rational approximation to A^{-1} based on an expansion in terms of the form $(A + \alpha I)^{-i}$.
 8 Because the matrix is shifted, its LU factorization might be a rather accurate representation of $A + \alpha I$.
 9 It is also more likely to be stable, in the sense defined in [6], in that its inverse does not have an
 10 extremely large norm. Instability of preconditioners in this sense is often the main cause of difficulty
 11 with incomplete LU factorization preconditioners of very ill-conditioned matrices. Thus, we can still
 12 manage to solve extremely ill-conditioned problems by exploiting incomplete factorizations of $A + \alpha I$,
 13 whereas an incomplete factorization for A would almost certainly result in failure.

14 Section 2 presents two algorithms for computing a rational preconditioner, which are illustrated on
 15 a simple example. Section 3 reports some error bounds for the conjugate gradient algorithm in the
 16 special case where the LU factorization of $A + \alpha I$ is computed exactly. A first error bound describes
 17 the behavior of the approximate solution at the beginning of the iteration process, and explains the
 18 fast decay of the error observed at the first steps in numerical experiments, even when using an
 19 incomplete factorization. A second error bound shows that the rational transformation improves the
 20 rate of convergence at the asymptotic regime, where, for example, a large accuracy is required.
 21 Finally, some numerical experiments on solving difficult real-world problems in structural mechanics
 are reported in Section 4.

23 2. Rational approximation preconditioning

Consider the linear system

$$Ax = b,$$

25 where A is a nonsingular square matrix of dimension n . A number of iterative methods approximate
 the solution $x = A^{-1}b$ to the above system, by a vector of the form

$$\tilde{x} = p(A)b,$$

27 where p is a certain polynomial. The approximation theory underlying these methods is to approxi-
 28 mate the rational function $s(\lambda) = 1/\lambda$ by a polynomial $p(\lambda)$ of degree d . The approximation is to be
 29 accurate on the (discrete) set of eigenvalues of A . However, preconditioners based on a polynomial
 approximation $p(A)$ of A^{-1} have their limitations, and as a result they are currently seldom used.

31 Rational approximations can be considered as an alternative. The first reaction to this possible
 32 approach is that the problem may not be well defined, since the best approximation to $1/\lambda$ by
 33 rational functions is $1/\lambda$ itself. A hint at a possible approach is to consider a similar situation that
 is naturally encountered when solving large eigenvalue problems. In shift-and-invert strategies [10],
 35 it is common to compute an eigenvalue λ_i by using a Krylov-subspace type method on the matrix
 $(A - \sigma I)^{-1}$, where σ is chosen to be close to the desired eigenvalue λ_i . More general rational Krylov

1 subspaces have also been used in [11] for eigenvalue computations. The goal is similar here since
 2 the inverse function is to be approximated by a rational function with a pole close to the origin.

3 2.1. Approximations to $1/\lambda$

4 We wish to obtain the best possible approximation to the function $s(\lambda) \equiv 1/\lambda$ from an expansion
 5 of the form

$$\frac{1}{\lambda} \simeq \eta_0 + \frac{\eta_1}{\lambda + \alpha} + \frac{\eta_2}{(\lambda + \alpha)^2} + \cdots + \frac{\eta_d}{(\lambda + \alpha)^d} + \cdots, \quad (1)$$

6 where $\alpha > 0$. We can use a Padé-type approximation to match both sides of (1), in the variable
 7 $t \equiv \lambda + \alpha$ to the approximation: we multiply both sides by t^d and require that the expansions in
 8 terms of t^j agree on both sides up to the highest possible degree. A little calculation yields the
 9 approximation

$$s(\lambda) \simeq \sum_{i=1}^d \frac{\alpha^{i-1}}{(\lambda + \alpha)^i}. \quad (2)$$

This can alternatively be verified by considering the following expansion:

$$\begin{aligned} \sum_{i=1}^d \frac{\alpha^{i-1}}{(\lambda + \alpha)^i} &= \frac{1}{\lambda + \alpha} \times \frac{1 - (\alpha/(\lambda + \alpha))^d}{1 - \alpha/(\lambda + \alpha)} \\ &= \frac{1}{\lambda} \left[1 - \left(\frac{\alpha}{\lambda + \alpha} \right)^d \right]. \end{aligned} \quad (3)$$

11 The relative error is given by

$$e(\lambda) = \left(\frac{\alpha}{\alpha + \lambda} \right)^d. \quad (4)$$

12 The function $\lambda s(\lambda)$ which gives an idea of the eigenvalues of the preconditioned matrix is illustrated
 13 in Fig. 1 for different values of α , with $d=3$ (left), and $d=6$ (right), using values of α ranging from
 14 0.025 to $\alpha = 0.375$ with increments of 0.05. The smaller α is, the closer the curve to the constant
 15 one. Thus, for the left plot, the highest curve (close to one) corresponds to $\alpha=0.025$, and the lowest
 16 one to $\alpha = 0.025 + 7 \times 0.05 = 0.375$. For the right plot ($d=6$) only the lowest 4 curves are shown.
 17 Notice that for large λ all curves are fairly accurate. For the smaller values of λ , the quality of the
 18 approximation is still excellent and stays close to one for small α .

19 2.2. Compounding ILU and shifting

Substitution of A for λ in (3) leads to an approximation to A^{-1} , given by

$$A^{-1} \simeq \sum_{i=1}^d \alpha^{i-1} (A + \alpha I)^{-i}. \quad (5)$$

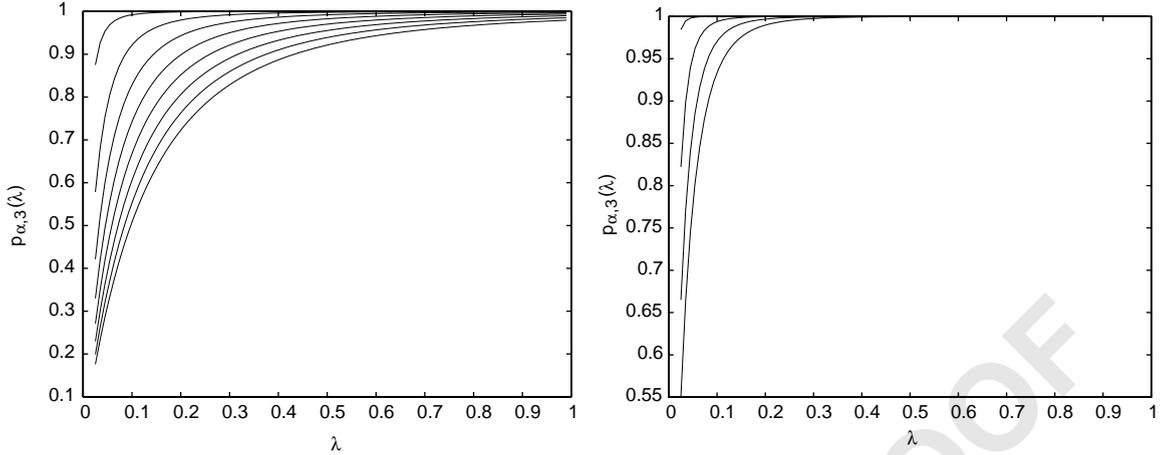


Fig. 1. The functions $\lambda s(\lambda)$ in the interval $[0.025, 1]$, for different values of α and $d = 3$ (left) and $d = 6$ (right).

- 1 The attraction of the above expansion is that it allows to ‘extrapolate’ an approximate LU factoriza-
 3 [14] of the matrix with a small α

$$A + \alpha I = L_\alpha U_\alpha + R_\alpha, \quad M_\alpha \equiv L_\alpha U_\alpha. \quad (6)$$

- 5 Then a rational preconditioning operation used for preconditioning the original matrix A can be
 5 defined by

$$M^{-1}v = \sum_{i=1}^d \alpha^{i-1} M_\alpha^{-i} v. \quad (7)$$

- 7 If the matrix M_α is symmetric, then the preconditioner M is also symmetric, and can be used by
 7 the conjugate gradient algorithm. An algorithm for computing $w = M^{-1}v$ is the following.

Algorithm 2.1. Symmetric rational preconditioning operation

1. $w := v$
2. *Do* $j = 1 : d - 1$
3. $w := v + \alpha M_\alpha^{-1} w$
4. *EndDo*
- 9 5. $w := M_\alpha^{-1} w$.

- 11 In order to relate this new technique with already known methods, consider the special case when
 11 an exact factorization is employed (i.e., $R_\alpha = 0$ in (6)) in the above algorithm. In this situation
 13 observe that $\alpha M_\alpha^{-1} = I - AM_\alpha^{-1}$ and as a result the preconditioning operation in line 3 of Algorithm
 13 2.1 can also be written as $w := v + (I - AM_\alpha^{-1})w$. This yields an alternative algorithm which can
 be viewed as a generalization of Algorithm 2.1.

1 **Algorithm 2.2.** Nonsymmetric rational preconditioning operation

1. $w := v$
2. *Do* $j = 1 : d - 1$
3. $w := v + (I - AM_x^{-1})w$
4. *EndDo*
5. $w := M_x^{-1}w$.

3 An approximate LU factorization M_x of $A + \alpha I$ can also be used in the same manner as before. However, the resulting preconditioner will no longer be symmetric in general, even if A is symmetric, but symmetry may be recovered in the usual way when M_x is available in the factored form $M_x = L_x L_x^T$, see, e.g., [14].

7 It is interesting to compare these two approaches, at least theoretically. Recall that the main idea of the rational approximation preconditioner is to approximate A^{-1} from expansions of an LU factorization of a near-by matrix, specifically $A + \alpha I$. We first consider the case when M_x represents an exact factorization of $A + \alpha I$. In that case, Algorithm 2.1 can be derived from the equality

$$A^{-1} = [(I - \alpha M_x^{-1})M_x]^{-1} = M_x^{-1}[I - \alpha M_x^{-1}]^{-1}$$

11 by applying a Neumann series expansion to approximately invert the bracketed term. In contrast, Algorithm 2.2 is based on approximating again A^{-1} , but through the identity

$$A^{-1} = [(I - (I - AM_x^{-1}))M_x]^{-1} = M_x^{-1}[I - (I - AM_x^{-1})]^{-1}$$

13 to which, again, a standard Neumann series expansion is applied to approximate the inverse of the second term of the right-hand side.

15 In the general case of an incomplete LU factorization, we have from (6)

$$I - AM_x^{-1} = (\alpha I - R_x)M_x^{-1}.$$

The preconditioning matrix of Algorithm 2.1 can be written as

$$M_x^{-1} \sum_{i=0}^{d-1} \alpha^i M_x^{-i} = M_x^{-1} (I - \alpha M_x^{-1})^{-1} (I - \alpha^d M_x^{-d}) \quad (8)$$

$$= (A - R_x)^{-1} (I - \alpha^d M_x^{-d}). \quad (9)$$

The preconditioning matrix of Algorithm 2.2 corresponds to

$$M_x^{-1} \sum_{i=0}^{d-1} (I - AM_x^{-1})^i = M_x^{-1} (AM_x^{-1})^{-1} (I - [I - AM_x^{-1}]^d) \quad (10)$$

$$= A^{-1} (I - [(\alpha I - R_x)M_x^{-1}]^d). \quad (11)$$

As is expected, Algorithms 2.1 and 2.2 coincide only when the factorization is exact ($R_x = 0$). Denote by P_1 the preconditioning matrix in (8) and by P_2 the preconditioning matrix defined by (10). The preconditioned matrices AP_i , $i = 1, 2$ are more relevant to understanding the quality of the

1 preconditioner. They are given by

$$AP_1 = (I - R_\alpha A^{-1})^{-1} (I - \alpha^d M_\alpha^{-d}), \quad (12)$$

$$AP_2 = I - [(\alpha I - R_\alpha) M_\alpha^{-1}]^d. \quad (13)$$

3 Since $\alpha I - R_\alpha = M_\alpha - A$, the preconditioned matrix AP_2 related to the second algorithm is fairly easy
 4 to analyze. Its eigenvalues are equal to $1 - (1 - \lambda)^d$, where λ represents a generic eigenvalue of the
 5 matrix AM_α^{-1} . In contrast, the eigenvalues of the preconditioned matrix AP_1 are not related in an
 6 easy way to those of $M_\alpha^{-1}A$, except in special situations, such as when $M_\alpha = A + \alpha I$.

7 We now compare formulas (12) and (13), in the general case when $R_\alpha \neq 0$. Consider first the
 8 situation when R_α is small, say much smaller than α . Then the residual matrix for the Neumann series
 9 preconditioner would be of the order of α^d in both cases, and so the accuracy of both algorithms
 10 is likely to be comparable. However, Algorithm 2.2 is more expensive than Algorithm 2.1 because
 11 of the extra matrix–vector product with A . An interesting case is when A is a dense matrix arising
 12 from electromagnetics and M_α is a sparse preconditioner to $A + \alpha I$. In this case, vector products with
 13 A are very expensive and Algorithm 2.1 is certainly advantageous.

14 Consider now the reverse situation when α is very small compared with R_α as measured by a
 15 certain norm. In this case, the error for Algorithm 2.1 is of the order of $\|R_\alpha\|$, while we obtain an
 16 order of $\|R_\alpha\|^d$ for Algorithm 2.2. Roughly speaking, this tells us that *when R_α is large relative to
 17 α , the use of Algorithm 2.1 will make little sense, because the potential gain from the expansion
 in α is annihilated by the large error in R_α .*

18 In summary, Algorithm 2.1 is likely to be competitive with the more traditional Algorithm 2.2
 19 only when α is large enough and R_α is small. In all cases when Algorithm 2.1 is employed, the ILU
 20 factorization of $A + \alpha I$ should be fairly accurate. If not, the resulting procedure can be ineffective,
 21 a fact that is confirmed by experiments. However, it is often not a problem to compute an accurate
 22 ILU factorization of $A + \alpha I$ provided an adequate shift α is applied. All these facts are illustrated
 23 on a simple example in the next sections.

2.3. A simple illustration

24 In order to illustrate the points of the above discussion we now consider a test example using MAT-
 25 LAB. The linear system considered simulates a problem which arises in the stream-function/vorticity
 26 formulation of the Navier–Stokes equations in two dimensions. As is well known [2,7] this formu-
 27 lation leads to a nonlinear equation of the form

$$F(\psi) = \Delta^2 \psi + Re[\psi_y(\Delta \psi)_x - \psi_x(\Delta \psi)_y].$$

28 Here Δ^2 represents the biharmonic operator, and u_x , u_y represents the partial derivatives of the
 29 function u with respect to x and y , respectively. The differential of the above system is represented
 30 by the linear operator:

$$\frac{DF}{D\psi} u \equiv \Delta^2 u + Re[(\Delta \psi)_x u_y - (\Delta \psi)_y u_x + \psi_y(\Delta u)_x - \psi_x(\Delta u)_y].$$

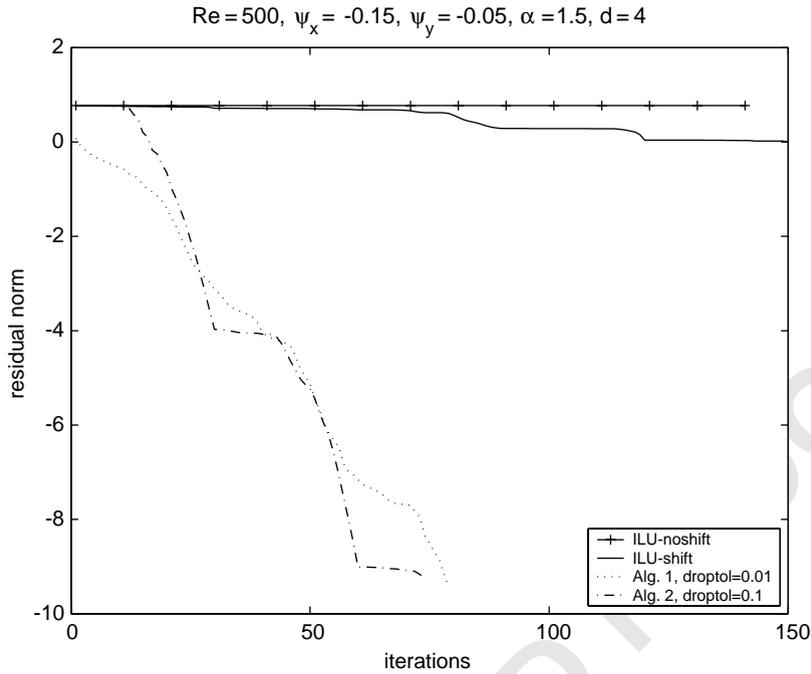


Fig. 2. Solution of the model stream function linear system of size 1225, solved by GMRES(30) preconditioned in 4 different ways.

1 preconditioner is now used in conjunction with Algorithm 2.1. Since the accuracy is rather poor, this
 2 method stagnates and the corresponding curve, which is not shown, is almost identical with the first
 3 (top) curve. Algorithm 2.2 on the other hand yields good convergence (dash-dotted curve). To verify
 4 the explanations given above, we performed another experiment with Algorithm 2.1 with a fairly
 5 accurate ILU factorization, one that is obtained from MATLAB with a drop tolerance $\text{droptol} = 0.01$.
 6 With this factorization, Algorithm 2.1 converges similarly to Algorithm 2.2 (dotted line). One might
 7 argue here that applying Algorithm 2.1 is more expensive. As was explained before, this is certainly
 8 not true when the matrix A is fairly dense. For this particular example, Algorithm 2.1 is in fact more
 9 advantageous than Algorithm 2.2. Indeed for $\text{droptol} = 0.1$ the total number of nonzeros for L and
 10 U together is $\text{nnz}(M) = 4761 + 3605 = 8636$ and this increases to $\text{nnz}(M) = 9303 + 8194 = 17,497$
 11 for the more stringent $\text{droptol} = 0.01$. The number of nonzero elements in the original matrix A is
 12 $\text{nnz}(A) = 15,229$. So each of the $d - 1$ sub-steps (represented by line 3 in both algorithms) costs
 13 $17,497 + 1225 = 18,722$ operations for Algorithm 2.1 and $15,229 + 8636 + 1225 = 24,840$ operations
 14 for Algorithm 2.2. For Algorithm 2.1, when $d = 4$, this is done 3 times and another solve with M is
 15 applied in line 5 leading to a total of $18,722 \times 3 + 17,497 = 73,663$ operations. For Algorithm 2.2,
 16 this count becomes $24,840 \times 3 + 8636 = 83,156$. Clearly, the advantage for Algorithm 2.1 improves
 17 as d increases. Another point to make here is that this problem is actually not a hard one to solve.
 18 A good ILU factorization is obtained for the nonshifted matrix when a small enough drop tolerance
 19 is used.

1 2.4. Inner–outer rational preconditioning

3 An appealing alternative to the above approach aims at extracting an optimal solution from the
 4 iterates of Algorithm 2.1. Instead of a preconditioning of the form (7) we may seek a preconditioned
 5 vector of the form

$$M^{-1}v = \sum_{i=1}^d \alpha_i M_\alpha^{-i} v,$$

6 where the scalars α_i are determined to minimize the residual norm $\|v - Aw\|_2$. In the case when M
 7 does not represent an accurate ILU factorization, we already know that such a sequence is not likely
 8 to be effective. It is therefore more general and effective to seek an optimal combination of iterates
 9 from Algorithm 2.2. This means that we seek a combination of the vectors of the preconditioned
 Krylov subspace

$$K_{\alpha,d} = \text{span}\{M_\alpha^{-1}Av, \dots, (M_\alpha^{-1}A)^d v\}.$$

The usual least-squares solution obtained by GMRES is calculated to minimize the residual norm.
 11 This method is nothing but an inner–outer GMRES iteration, in which the inner solve is itself
 12 a GMRES iteration using a Krylov subspace of dimension d . In this case, the preconditioner is
 13 modified at each step, hence a flexible accelerator such as FGMRES [12] must be used. Numerical
 14 results indicate that this is more costly but often more effective than the simple expansion (7) (see
 15 Fig. 5).

16 2.5. A multiscale-type procedure using different shifts

17 It is often observed that after a certain number of steps the convergence of GMRES slows down
 18 considerably, sometimes to the point of stagnating. This usually means that certain modes are not
 19 captured by the iterative process. Assume that the incomplete factorization is exact and consider

$$A + \alpha I = L_\alpha U_\alpha.$$

20 According to (9) with $R_\alpha = 0$, the preconditioning matrix which is defined by Algorithm 2.1 is equal
 21 to

$$M^{-1} = A^{-1}(I - \alpha^d(A + \alpha I)^{-d})$$

and so the residual matrix, which is

$$I - AM^{-1} = \alpha^d(A + \alpha I)^{-d}$$

22 has eigenvalues:

$$\rho_i = \left(\frac{\alpha}{\lambda_i + \alpha} \right)^d,$$

23 where λ_i is an arbitrary eigenvalue of A . These are the same functions as the errors in (4) (hence Fig.
 24 1 plots also the function $1 - \rho_i$). It is clear that for large α those residual components associated with
 25 eigenvalues λ_i that are close to zero will not be reduced much. Notice that any eigenvalue outside
 26 the disk $\bar{D}(-\alpha, \alpha)$ centered at $-\alpha$ and of radius α , will be damped, i.e., it will be transformed into

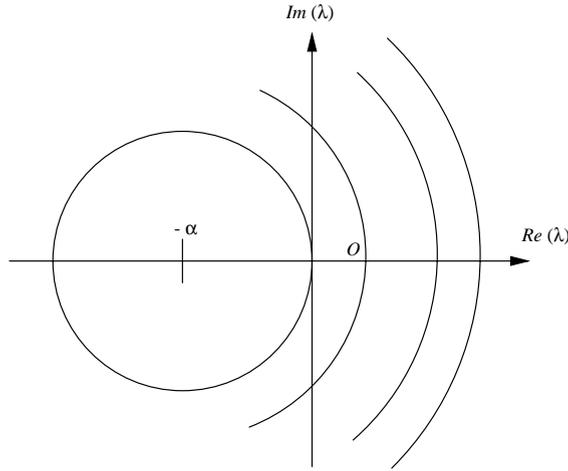


Fig. 3. Damping region for preconditioner is outside the disk.

1 an eigenvalue smaller than one. Eigenvalues inside the disk can cause serious difficulties since they
 2 can be amplified and become very large if d is large.
 3 Assuming that there are no eigenvalues inside the disk (as is the case for positive definite ma-
 4 trices), all damping ratios will be less than one. The farther away λ_i is from the center $-\alpha$ the
 5 smaller will be the damping ratio. Those eigenvalues close to the circle of center $-\alpha$ and radius
 6 α will have a damping ratio close to one. The concentric arcs in Fig. 3 show the lines where the
 7 eigenvalues have the same damping factor ρ . If α is not changed during the iteration process, then
 8 the eigenvectors of the residual which have small damping ratio (corresponding to large eigen-
 9 values) will be eliminated quickly. Those with damping ratios close to one (for example those close
 10 to the origin) are likely to change very little. What would be ideal is to have a procedure that does
 11 not disturb those small residual components achieved in earlier steps, but that reduces those closer
 12 to the origin further. This can be easily done by reducing α , say, at the occasion of a restart in
 13 the GMRES(m) algorithm. Experiments do indeed show that this principle works: similar to other
 14 multiscale methods, it is much better to work on different parts of the spectrum—using different
 15 stages—rather than using a preconditioner based on a single α .

2.6. Tests with inner–outer techniques and variable shift

17 We now go back to the example of Section 2.3 to illustrate the techniques of the previous two
 18 sections. For the variable shift strategy of Section 2.5, we do not use a different factorization for
 19 each different α , since this would be prohibitive. Instead, the same LU factorization as for the other
 20 methods is used, i.e., the one obtained with the shift $\alpha = 1.5$ and using the same drop tolerance
 21 of 0.01. The variable shift strategy is used in conjunction with Algorithm 2.2. The variation of α
 22 occurs in line 3 which is replaced by $w := v + (I - (A + \sigma I)M_\alpha^{-1})w$. The resulting preconditioner
 23 becomes

$$(A + \sigma I)^{-1}(I - [([\alpha - \sigma]I - R_\alpha)M_\alpha^{-1}]^d)$$

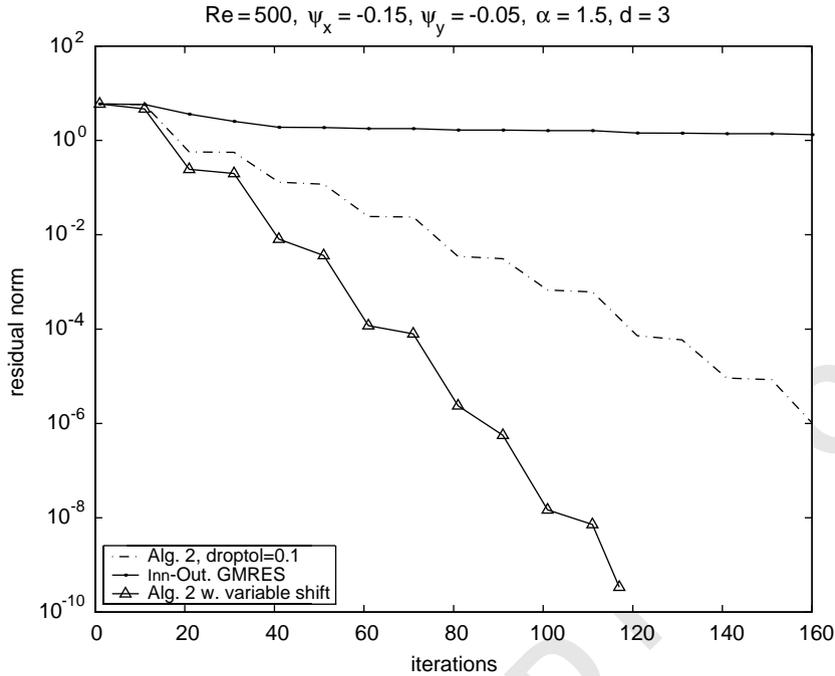


Fig. 4. Solution of the model stream function linear system with $Re = 500$.

1 instead of (see Eq. (11))

$$A^{-1}(I - [(\alpha I - R_x)M_x^{-1}]^d).$$

3 Hence α has been changed to $\alpha - \sigma$ in the right parentheses, with $0 < \sigma < \alpha$, thus (partially) de-
 5 creasing α as proposed in Section 2.5. Regarding the shifting strategy, the initial shift σ is set to
 7 $\sigma = 0.05$ then at each FGMRES outer iteration σ is increased by $\delta\sigma \equiv 0.01$. Other strategies for
 9 shifting have been tested but did not yield better results.

7 Fig. 4 compares three different methods: Algorithm 2.2, inner–outer GMRES of Section 2.4 and the
 9 variable- α procedure of Section 2.5. In order to have shorter restarts, the Krylov subspace dimension
 11 was reduced from $m = 30$ to 20. All other parameters remain the same as in the example of Section
 13 2.3. In particular, the shift α is again 1.5, and the drop tolerance for `iluinc` is `droptol=0.1`, yielding
 15 the same ILU factorization. The degree d has been set to $d = 3$. This is also the subspace dimension
 17 for the inner GMRES iteration in inner–outer GMRES algorithm. In this particular case, the inner–
 outer GMRES iteration does not converge. Algorithm 2.2 (with fixed shift) behaves similarly to
 the case with $m = 30$ shown in Section 2.3. The interesting observation is that the variable shift
 Algorithm does much better than the other 2 algorithms. Note that this improvement is achieved with
 a small change to Algorithm 2.2, since in most iterations the shift σ is rather small ! Observe also
 that this method uses about the same number of operations per iteration as Algorithm 2.2 (without
 shift) and fewer than the inner–outer GMRES method. Although only iteration counts are shown,
 the plots in Figs. 4 and 5 do provide some idea of the operations cost in some cases.

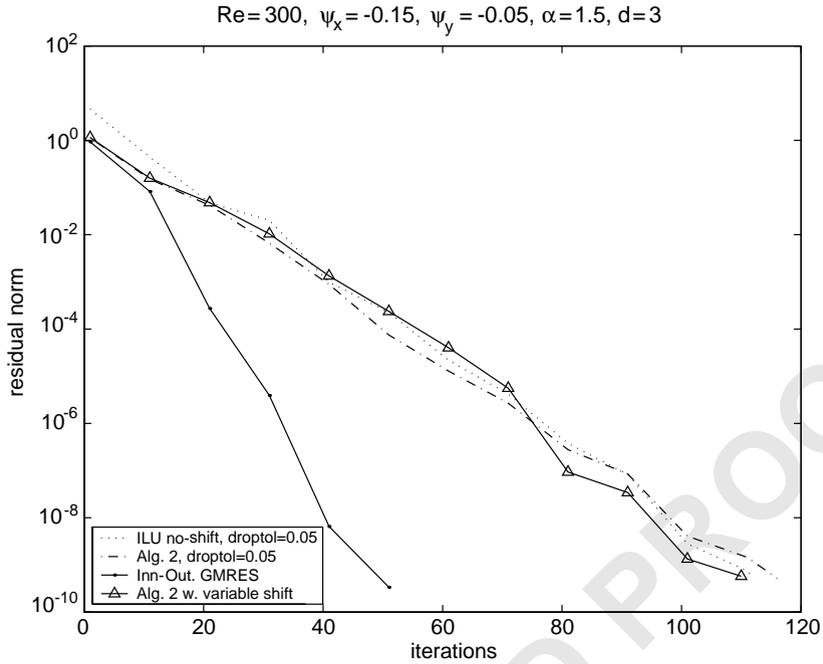


Fig. 5. Solution of the model stream function linear system with $Re = 300$.

1 As was mentioned above, the inner–outer GMRES iteration did not perform too well in this
 2 example. In our experience we found that this happened occasionally but that this occurrence is by
 3 no means general. For example, in the next test we changed the Reynolds number to $Re = 300$ and
 4 computed the ILU factorization with a better accuracy using $\text{droptol} = 0.05$. In that case, we had
 5 $\text{cond}_{\text{est}}(\text{LU}) = 3.1 \times 10^5$ without shift and $\text{cond}_{\text{est}}(\text{LU}) = 138$ with shift. The ILU factorization for
 6 the matrix A (without shift) is accurate and fairly stable—so the corresponding ILU preconditioning
 7 works well in this case. The corresponding run is shown in a dotted line in Fig. 5. As is expected
 8 in this case not much can be gained from the rational approximation preconditioner. In particular,
 9 the variable shift technique is no longer competitive with the inner–outer GMRES technique which
 10 does quite well.

11 3. Analysis in the case of an exact factorization

12 This section gives some error bounds for the CG algorithm applied to the solution of the pre-
 13 conditioned system

$$Bx := M^{-1}Ax = M^{-1}b. \quad (15)$$

14 The rational transformation $B = r(A)$ makes the spectrum of B clustered around 1, and the usual
 15 error bound for the CG algorithm can be improved, both for the first iterations and asymptotically.

16 Here it is assumed that the LU factorization is exact, i.e., $R_x = 0$ in (6). As was already seen,
 17 this means that Algorithms 2.1 and 2.2 coincide. The case of an approximate LU factorization is

1 more difficult to analyze, though one can expect the behavior to be similar if the perturbation of the
 2 exact factorization remains small. As was mentioned in Section 2.2 and observed in the numerical
 3 experiments in Sections 2.3 and 2.6, the ILU factorization of $A + \alpha I$ should be fairly accurate in
 order to obtain an effective procedure.

5 The matrix A is assumed to be symmetric positive definite, with increasingly ordered
 eigenvalues λ_i ,

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = 1,$$

7 associated to normalized eigenvectors v_i , $i = 1, \dots, n$. The preconditioning operation corresponding
 to (3) is modified by a constant β into

$$r(\lambda) = \frac{1}{\beta} \left[1 - \left(\frac{\alpha}{\lambda + \alpha} \right)^d \right], \quad \beta = 1 - \left(\frac{\alpha}{1 + \alpha} \right)^d, \quad (16)$$

9 where $r(0) = 0$, $r(1) = 1$, and the eigenvalues of $B := r(A)$ are

$$\mu_i = r(\lambda_i), \quad 1 \leq i \leq n. \quad (17)$$

11 The choice of β ensures that $\mu_n = r(\lambda_n) = 1$. As the matrix A is assumed to be ill-conditioned, the
 eigenvalue λ_1 is very close to 0, and the shift α can be chosen greater than λ_1 , but still noticeably
 13 smaller than 1. Hence $\beta \simeq 1$. The matrix B remains symmetric positive definite, and has the same
 eigenvectors v_i as A .

Let $x = A^{-1}b$ be the exact solution. We denote by x_m the approximate solution obtained at the
 15 m th step of the CG algorithm applied to the matrix A , and by y_m the approximate solution obtained
 with the matrix B , starting with $y_0 = x_0$. A well-known optimality property of the CG algorithm,
 17 states that

$$\|y_m - x\|_B^2 = \min_{p \in \mathbb{P}_m, p(0)=1} \|p(B)(x_0 - x)\|_B^2, \quad (18)$$

19 where \mathbb{P}_m is the set of polynomials of degree at most m , and $\|\cdot\|_B$ is the B -norm defined by
 $\|z\|_B^2 = (Bz, z)$. The minimizer p_m of (18) and the vector y_m are related by $y_m - x = p_m(B)(x_0 - x)$.

21 The classical error bound for the CG algorithm applied to the solution of the system $Ax = b$ is
 given by

$$\|x_m - x\|_A^2 \leq 4 \left(\frac{1 - \sqrt{\lambda_1}}{1 + \sqrt{\lambda_1}} \right)^{2m} \|x_0 - x\|_A^2, \quad (19)$$

whereas the error for the preconditioned system (15) is

$$\|y_m - x\|_B^2 \leq 4 \left(\frac{1 - \sqrt{\mu_1}}{1 + \sqrt{\mu_1}} \right)^{2m} \|x_0 - x\|_B^2. \quad (20)$$

23 3.1. Error bounds for the first iterations

For given $k \geq 0$ and $m \geq 1$, let

$$p(t) = \frac{T_m[1 + 2(\mu_{k+1} - t)/(1 - \mu_{k+1})]}{T_m[1 + 2\mu_{k+1}/(1 - \mu_{k+1})]} \quad (21)$$

1 be the scaled Chebyshev polynomial that is small in the interval $[\mu_{k+1}, 1]$ (recall that $\mu_n = 1$). It satisfies $p(0) = 1$ and the following bounds, see e.g., [14]:

$$p^2(t) \leq 1, \quad \forall t \in [0, 1], \quad (22)$$

$$p^2(t) \leq \zeta := 4 \left(\frac{1 - \sqrt{\mu_{k+1}}}{1 + \sqrt{\mu_{k+1}}} \right)^{2m} \quad \forall t \in [\mu_{k+1}, 1]. \quad (23)$$

3 We have

$$\|p(B)(x_0 - x)\|_B^2 = \sum_{i=1}^n p^2(\mu_i) e_i^2 \mu_i, \quad e_i^2 \mu_i \geq 0,$$

5 where $e_i = (x_0 - x, v_i)$ is the initial error projected on the eigenvector v_i . The initial error $\|x_0 - x\|_B^2$ can be split into two parts as follows:

$$\|x_0 - x\|_B^2 = S_1 + S_2,$$

$$S_1 = \sum_{i=1}^k e_i^2 \mu_i, \quad S_2 = \sum_{i=k+1}^n e_i^2 \mu_i.$$

Then, using (22) and (23), we obtain

$$\|p(B)(x_0 - x)\|_B^2 = \sum_{i=1}^k p^2(\mu_i) e_i^2 \mu_i + \sum_{i=k+1}^n p^2(\mu_i) e_i^2 \mu_i \leq S_1 + \zeta S_2,$$

7 from which it follows that

$$\|y_m - x\|_B^2 \leq S_1 + 4 \left(\frac{1 - \sqrt{\mu_{k+1}}}{1 + \sqrt{\mu_{k+1}}} \right)^{2m} \|x_0 - x\|_B^2. \quad (24)$$

9 Observe that the special choice $k = 0$ leads to $S_1 = 0$ which yields inequality (20) as a particular case.

11 Up to now, we have not exploited the clustering of the spectrum of $B = r(A)$ around 1. A consequence of this clustering is that $S_1 = S_1(k)$ may be small while at the same time μ_{k+1} is close to 1. We have $\mu_{k+1} = r(\lambda_{k+1})$ with

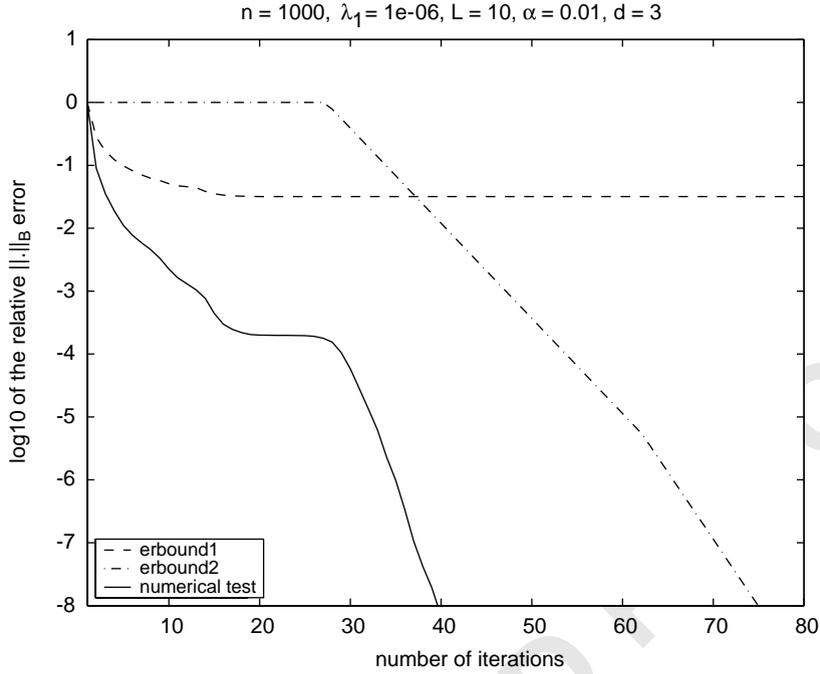
$$\lambda_{k+1} = c\alpha$$

13 for a certain positive number c . We must keep in mind that we want to use an α larger than λ_1 but smaller than 1, which gives the possible range for c . We now define L to be the smallest scalar independent of α and k , for which

$$\sum_{i=1}^k e_i^2 \leq \lambda_k L \|x_0 - x\|^2 \quad \forall k = 1, 2, \dots, n. \quad (25)$$

17 Recalling that $\|x_0 - x\|^2 = \sum_{i=1}^n e_i^2$, note that the linear function $l(\lambda) \equiv L\lambda$ can be viewed as an upper bound for the density function, a step function whose value at each eigenvalue λ_j is defined by

$$\phi(\lambda_j) = \frac{\sum_{i=1}^j e_i^2}{\sum_{i=1}^n e_i^2}.$$

Fig. 6. Error bounds $\|y_m - x\|_B$.

- 1 For example, if the matrix A has a uniform spectrum $\lambda_i = i/n$ and if $e_i^2 = 1$ for all i , then (25) holds
 2 for $L = 1$. A large value of L corresponds to a clustering of the spectrum of A near 0, or to an
 3 initial error $x_0 - x$ essentially concentrated on the eigenspace associated with the smallest eigenvalues.
 4 Then we have the following error bound, which explains the fast decay of the error observed at the
 5 beginning of the iteration process (see Figs. 6 and 7).

7 **Proposition 3.1.** For $\alpha > 0$, let $c > 0$ be chosen such that $c\alpha = \lambda_{k+1}$ is an eigenvalue λ_{k+1} of the
 8 matrix A . Then

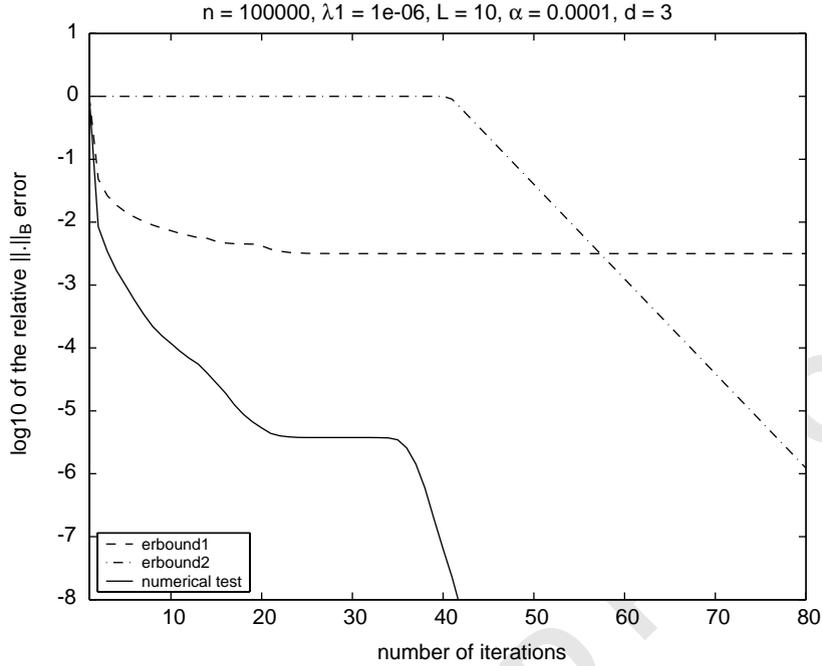
$$\|y_m - x\|_B^2 \leq \left[c\alpha L + 4 \left(\frac{1}{\sqrt{(c+1)^d} + \sqrt{(c+1)^d - 1}} \right)^{4m} \right] \|x_0 - x\|^2. \quad (26)$$

Proof. Since $\mu_i \leq 1$ and by assumption (25), we have

$$S_1 = \sum_{i=1}^k e_i^2 \mu_i \leq \sum_{i=1}^k e_i^2 \leq \lambda_k L \|x_0 - x\|^2.$$

- 9 It follows from (24), $\lambda_k \leq \lambda_{k+1} = c\alpha$, $\|\cdot\|_B \leq \|\cdot\|$ and $\mu_{k+1} = r(\lambda_{k+1})$ that

$$\|y_m - x\|_B^2 \leq \left[c\alpha L + 4 \left(\frac{1 - \sqrt{r(c\alpha)}}{1 + \sqrt{r(c\alpha)}} \right)^{2m} \right] \|x_0 - x\|^2.$$

Fig. 7. Error bounds $\|y_m - x\|_B$.

1 For $0 < r < 1$, and since $\beta r < r$, we have

$$\frac{1 - \sqrt{r}}{1 + \sqrt{r}} \leq \frac{1 - \sqrt{\beta r}}{1 + \sqrt{\beta r}} = \left(\frac{\sqrt{1 - \beta r}}{1 + \sqrt{\beta r}} \right)^2,$$

then, using

$$\beta r(c\alpha) = 1 - 1/(c+1)^d,$$

3 we obtain

$$\begin{aligned} \frac{1 - \sqrt{r(c\alpha)}}{1 + \sqrt{r(c\alpha)}} &\leq \left(\frac{\sqrt{1/(c+1)^d}}{1 + \sqrt{1 - 1/(c+1)^d}} \right)^2 \\ &= \left(\frac{1}{\sqrt{(c+1)^d} + \sqrt{(c+1)^d - 1}} \right)^2. \end{aligned} \quad (27)$$

which completes the proof. \square

5 We can notice that if we use the polynomial

$$p(t) = \frac{T_j[1 + 2(\mu_{k+1} - t)/(1 - \mu_{k+1})]}{T_j[1 + 2\frac{\mu_{k+1}}{1 - \mu_{k+1}}]} \times \frac{T_{m-j}[1 + 2(\mu_1 - t)/(1 - \mu_1)]}{T_{m-j}[1 + 2\mu_1/(1 - \mu_1)]},$$

1 instead of (21), then the following bound is obtained for $1 \leq j \leq m$:

$$\|y_m - x\|_B^2 \leq 4 \left[c\alpha L + 4 \left(\frac{1}{\sqrt{(c+1)^d} + \sqrt{(c+1)^d - 1}} \right)^{4j} \right] \left(\frac{1 - \sqrt{\mu_1}}{1 + \sqrt{\mu_1}} \right)^{2(m-j)} \|x_0 - x\|^2.$$

However, for large m , the following section shows that a much stronger bound can be obtained.

3 3.2. Asymptotic error bounds

5 A strategy based on an idea described in [13], which takes advantage of the clustering of the
6 spectrum of $B = r(A)$ around one, is now used for obtaining asymptotic error bounds. Instead of
7 standard Chebyshev polynomials that are small in the interval $[\mu_1, \mu_n]$ we will use the following
8 modified polynomial:

$$C_m(t) = \prod_{i=1}^k \left(\frac{\mu_i - t}{\mu_i} \right) \times \frac{T_{m-k}[1 + 2(\mu_{k+1} - t)/(\mu_n - \mu_{k+1})]}{T_{m-k}[1 + 2(\mu_{k+1})/(\mu_n - \mu_{k+1})]}.$$

9 This consists of two parts. The first is a product term which takes the value zero for the first k
10 smallest eigenvalues μ_1, \dots, μ_k . The second is a standard scaled Chebyshev polynomial which is
11 small in the interval $[\mu_{k+1}, \mu_n]$. Note that C_m is of degree m and that $C_m(0) = 1$. Since $C_m(\mu_i) = 0$
for $i = 1, \dots, k$, the maximum of C_m on the spectrum of B is

$$\max_{\mu_j \in \Lambda(B)} |C_m(\mu_j)| \leq \max_{j=k+1, \dots, n} \prod_{i=1}^k \left| \frac{\mu_i - \mu_j}{\mu_i} \right| \times \frac{1}{T_{m-k}[1 + 2(\mu_{k+1})/(\mu_n - \mu_{k+1})]}. \quad (28)$$

13 **Proposition 3.2.** For $\alpha > 0$, let $c > 0$ be chosen such that $c\alpha = \lambda_{k+1}$ is an eigenvalue λ_{k+1} of the
matrix A . Then, for $m \geq k$, we have

$$\|y_m - x\|_B \leq \frac{2\kappa^k(\alpha)}{[\sqrt{(c+1)^d} + \sqrt{(c+1)^d - 1}]^{2(m-k)}} \|x_0 - x\|_B, \quad (29)$$

15 where

$$\kappa(\alpha) = \frac{1 - ((\lambda_1 + \alpha)/(1 + \alpha))^d}{((\lambda_1 + \alpha)/\alpha)^d - 1}. \quad (30)$$

Proof. Consider first the product term

$$\max_{j=k+1, \dots, n} \prod_{i=1}^k \left| \frac{\mu_i - \mu_j}{\mu_i} \right| = \prod_{i=1}^k \left| \frac{\mu_i - \mu_n}{\mu_i} \right| \leq \left[\frac{\mu_n - \mu_1}{\mu_1} \right]^k.$$

17 Recall that $\mu_n = 1$. We denote by $\kappa(\alpha)$ the term $(1 - \mu_1)/\mu_1$ inside the parentheses:

$$\kappa(\alpha) = \frac{1 - \frac{1}{\beta}(1 - \alpha^d/(\lambda_1 + \alpha)^d)}{\frac{1}{\beta}(1 - \alpha^d/(\lambda_1 + \alpha)^d)}$$

$$\begin{aligned}
&= \frac{1 - \alpha^d/(1 + \alpha)^d - 1 + \alpha^d/(\lambda_1 + \alpha)^d}{1 - \alpha^d/(\lambda_1 + \alpha)^d} \\
&= \frac{\alpha^d/(\lambda_1 + \alpha)^d - \alpha^d/(1 + \alpha)^d}{1 - \alpha^d/(\lambda_1 + \alpha)^d}.
\end{aligned}$$

1 Multiplying numerator and denominator by $(\lambda_1 + \alpha)^d/\alpha^d$ yields (30). Next, the second term is bounded by

$$\frac{1}{T_{m-k}[1 + 2(\mu_{k+1})/(1 - \mu_{k+1})]} \leq 2 \left(\frac{1 - \sqrt{\mu_{k+1}}}{1 + \sqrt{\mu_{k+1}}} \right)^{m-k}$$

3 and, using (27) once more, we obtain for $\mu_{k+1} = r(c\alpha)$

$$\frac{1}{T_{m-k}[1 + 2(\mu_{k+1})/(\mu_n - \mu_{k+1})]} \leq 2 \left(\frac{1}{\sqrt{(c+1)^d} + \sqrt{(c+1)^d - 1}} \right)^{2(m-k)}$$

which completes the proof. \square

5 A fair comparison would be between one step of the rational preconditioner versus d steps of
 7 the standard preconditioned conjugate gradient applied with some accurate ILU preconditioner. This
 is because applying $r(A)$ uses d solves with the LU factorization—which is likely to dominate the
 cost. For these d steps the convergence factor as inferred from the standard bound is given by

$$\tilde{\rho} = \left(\frac{1 - \sqrt{\lambda_1}}{1 + \sqrt{\lambda_1}} \right)^d$$

9 which is to be compared with the asymptotic convergence factor,

$$\rho(c) \leq \left(\frac{1}{\sqrt{(c+1)^d} + \sqrt{(c+1)^d - 1}} \right)^2.$$

The above asymptotic argument ignores the potentially large constant in the numerator of (29).

11 However, it gives a rough comparison of the situation at the asymptotic regime where, for example,
 a large accuracy is required.

13 The estimated error bounds resulting from Propositions 3.1 and 3.2 are illustrated by the curves
 erbound1 and erbound2 in Figs. 6 and 7 for two sets of parameters: $d = 3$, $\lambda_1 = 10^{-6}$, $L = 10$,
 15 $n = 10^3$, $\alpha = 10^{-2}$ in Fig. 6, and $d = 3$, $\lambda_1 = 10^{-10}$, $L = 10$, $n = 10^5$, $\alpha = 10^{-4}$ in Fig. 7. The other
 eigenvalues were chosen as follows:

$$\lambda_k = \max(\lambda_1 + (k - 1)/(nL), (2k - n)/n), \quad 2 \leq k \leq n.$$

17 The two curves erbound1 and erbound2 complement each other since the first estimate is better
 than the second one for the earlier iterates. They are compared to numerical results obtained from
 19 the CG algorithm preconditioned with Algorithm 2.1. The matrix A is of the form $A = Q^T \Lambda Q$ where
 Q is a randomly chosen sparse orthogonal matrix and Λ is diagonal with $\Lambda_{kk} = \lambda_k$. The factorization
 21 was not completely performed: we used the MATLAB command $U = \text{cholinc}(A + \alpha * I, \text{droptol})$
 with $\text{droptol} = \alpha^2$.

Table 1
Description of test problems

Name	n	n_z	Dominance (%)	Symmetry	Matrix source
ELTCOQUE	38,002	949,452	0.6	Yes	Shell modeling
MCHLNF	49,800	4,136,484	5	Yes	Tire design
MCHLNE	49,800	4,136,580	4.6	No	Tire design

1 4. Rational acceleration for realistic problems

3 This section reports on a few numerical experiments (on a DEC Alpha) with the rational pre-
 4 conditioning techniques for solving difficult real-world problems in structural mechanics. For these
 5 problems, a straightforward application of standard preconditioning techniques, such as an incomplete
 6 LU factorization, fails due to their instability. Diagonal shifting and large fill-in may be needed to
 7 achieve convergence as reported in [16] for tire design problems. However, choosing the best shift
 8 value can be time consuming, and the preconditioning becomes quite expensive to apply in the case
 9 of large fill-in. We attempt to show how rational preconditioning (Algorithm 2.1 or 2.2) can handle
 10 these difficulties. Since the purpose of the experiments shown here is to solve realistic applications
 11 problems more efficiently using a form of rational approximation, we present the experimental re-
 12 sults for either Algorithm 2.1 or 2.2 depending on which algorithm we found to perform better for
 13 a particular problem. For a comparison of the two algorithms, see Section 2.3.

13 4.1. Test problems and the components of the iterative solution

15 For the experiments, we have selected a few linear systems arising in shell modeling and tire
 16 design. Table 1 gives some information about the problems. The matrix ELTCOQUE comes from the
 17 discretization of thin shells, using DKT12 elements (Discrete Kirchoff Triangle with 12 ddl, [1]), and
 18 was provided by CADOE S.A. The matrices MCHLNF and MCHLNE come from the discretization
 19 of nonlinear static equilibrium equations in tire problems, and were provided by Michelin Americas
 20 Research and Development Corporation. Columns denoted by n and n_z give the numbers of rows
 21 and nonzero entries in the matrices, respectively. Column Dominance shows the ratio of diagonally
 22 dominant rows to the matrix size. This number gives a good indication of the difficulty of the
 23 corresponding linear system. The matrices have a small percent of the diagonally dominant rows,
 24 and thus we can expect the linear systems to be quite difficult. All the matrices in Table 1 are
 25 *structurally* symmetric. The symmetry in value is indicated in Column Symmetry.

26 Restarted GMRES was used as the accelerator. Specifically, the FGMRES(k) variant, which al-
 27 lows variable preconditioning [14], was employed in cases when the preconditioner changes during
 28 iteration. Deflated GMRES(k) [3,9] was used in cases of stagnation. In deflated GMRES(k) the
 29 eigenvectors corresponding to a few smallest eigenvalues are added to the Krylov subspace to pre-
 30 vent stalling of the GMRES(k) convergence. For both FGMRES(k) and deflated GMRES(k), the
 31 Krylov subspace dimension is equal to 54 and includes four injected eigenvectors in the case of
 deflated GMRES(k). We took a random initial guess with the right-hand side constructed such that
 the solution is the vector of all ones. For these problems, the effects of the artificial right-hand side

1 and the one coming from the application have been studied in [16]. It has been observed that the
 2 artificial right-hand side does not make these problems easier and that both right-hand side types
 3 result in a similar convergence pattern.

4 Rational acceleration is applied to the factorization produced by the algebraic recursive multilevel
 5 solver (ARMS) [15]. This choice of the preconditioner is motivated by the versatility of ARMS and
 6 its ability to solve efficiently the structural mechanics problems. ARMS is an algebraic multigrid-like
 7 algorithm that requires no underlying set of grids for defining prolongation and restriction operators.
 ARMS works by reordering the matrix in the block form

$$\begin{pmatrix} B & F \\ E & C \end{pmatrix},$$

8 in which B is diagonal or block-diagonal with small blocks. The above matrix is then approximately
 9 block-factored as

$$\begin{pmatrix} B & F \\ E & C \end{pmatrix} \approx \begin{pmatrix} L & 0 \\ G & I \end{pmatrix} \begin{pmatrix} U & W \\ 0 & S \end{pmatrix}$$

10 using again dropping strategies. Then the reordering and factorization were repeated recursively on
 11 the Schur complement matrix S , for a small number of levels. At the last level the matrix S is factored
 12 using again a standard ILUT or ILUTP factorization. Both the construction of the preconditioner
 13 and the forward-backward solutions in ARMS are recursive. In addition ARMS allows inter-level
 14 iterations (referred to as W-cycles in the multigrid literature), though these tend to be fairly expensive
 15 if the number of levels is high.

16 A particular instance of the ARMS preconditioner as well as the ARMS performance for a given
 17 iterative algorithm are controlled by several parameters, such as the block size and number of levels
 18 specifying the block and level preconditioner structures, respectively. We allow no inner iterations
 19 in the levels of ARMS to reduce the time of the preconditioning operation. Varying the number of
 20 ARMS levels from 2 to 5 did not affect the preconditioner performance, but fewer levels make the
 21 preconditioner construction less expensive. Thus, the number of levels was chosen to be equal to
 22 two. Our experience shows that taking small blocks (of size 3 or 10) instead of larger blocks (say,
 23 of size 100) often yields better overall performance.

24 Filtering small (less the 10^{-3}) off-diagonal entries in the matrix from which the preconditioners
 25 are built speeds up their construction since fewer nonzero entries remain in the original and pre-
 26 conditioner matrices. We have observed that in the given problem types after such a filtering process,
 27 the majority of (weakly) diagonally dominant rows have *all* their off-diagonal entries dropped. The
 28 corresponding rows constitute an independent set, which we call the *trivial independent set*. These
 29 rows are properly permuted by setting the independent set tolerance in ARMS to 1, see [15].

30 4.2. Rational acceleration and the accuracy of preconditioning

31 With a rational acceleration (Algorithm 2.1), a less accurate preconditioning matrix (i.e., with a
 32 small fill-in) may be sufficient to achieve a good convergence. For the MCHLNF and MCHLNE
 33 problems, Table 2 shows the results of the three runs of an experiment in which the accelera-
 34 tion degree was increased in each run while the amount of preconditioner fill-in was halved. The
 35

Table 2
Dependence of the execution times on the degree of rational acceleration and the preconditioner accuracy

Name	(Degree, Fill-in)	n_z	Construction	Solution	Iterations
MCHLNF	(2,240)	22,161,175	1327.12	1988.43	564
	(3,120)	11,301,437	726.81	1417.35	465
	(4,60)	5,710,647	410.15	1412.13	541
MCHLNE	(2,240)	22,161,330	1286.47	2256.35	626
	(3,120)	11,286,156	696.18	1575.95	508
	(4,60)	5,621,009	381.48	1406.52	550

1 total number of nonzero elements in the preconditioning matrix is shown in column 2. Columns
 2 Construction and Solution give the preconditioner construction and solution times (in seconds),
 3 respectively. The number of *outer* deflated GMRES(54) iterations is stated in column 5. The shift
 4 value α and the ARMS dropping tolerance have been kept constant and equal to 0.8 and 0.0, re-
 5 spectively. The reduction of 10^6 in the residual norm has been achieved by deflated GMRES(54).

6 As expected, the preconditioner construction time is almost proportional to the amount of fill-in
 7 and affects significantly the total execution time. The preconditioning application cost itself increases
 8 when the degree d of approximation grows, but decreases when the amount of fill-in is reduced.
 9 Similarly, increasing d reduces the number of outer iterations, while reducing the fill-in augments
 10 the number of outer iterations. These opposite tendencies finally result in a reduction of the solution
 11 time. In our experiments, we have observed that, for this problem, halving the amount of fill further
 12 makes the preconditioner very inaccurate, and the solution time grows along with the added cost of
 13 increased degree d . Hence, a reasonable increase of the degree and reduction in the fill-in reduces
 14 both the construction and solution time, that is, in this example, the rational acceleration saves time
 15 and memory.

4.3. Shift and degree selection

17 Here, we show the dependences of the convergence rate and the stability of the preconditioner on α
 18 and outline an automatic process of arriving at an appropriate shift value. In [4], a strong correlation
 19 between stability of the preconditioner and the size of $\mathcal{E} = \log(\|(LU)^{-1}\|_{\text{inf}})$ was shown and was
 20 suggested as a practical means of evaluating the quality of a preconditioner. We can inexpensively
 21 compute \mathcal{E}_α as

$$\mathcal{E}_\alpha = \log(\|(LU)^{-1}e\|_1),$$

where e is a vector of all ones and LU is the incomplete LU factorization of $A + \alpha I$.

23 For the problem ELTCOQUE, the amount of fill-in was equal to 30 and the dropping tolerance
 24 equal to 0 in the preconditioner construction. Without rational preconditioning and without shift
 25 ($\alpha = 0$), there was no reduction of the residual norm. Fig. 8 shows the convergence curves for
 26 different choices of α with the lowest degree rational approximation (i.e., its degree $d = 1$). When α
 27 is quite large (solid line) the convergence of flexible GMRES(54) is slow although the incomplete
 28 LU factors are stable ($\mathcal{E}_\alpha = 0.29$). It is possible to start with some large α (say, 0.8) and gradually
 29 decrease it as long as the indicator \mathcal{E}_α stays small. Changing α dynamically requires modifying the

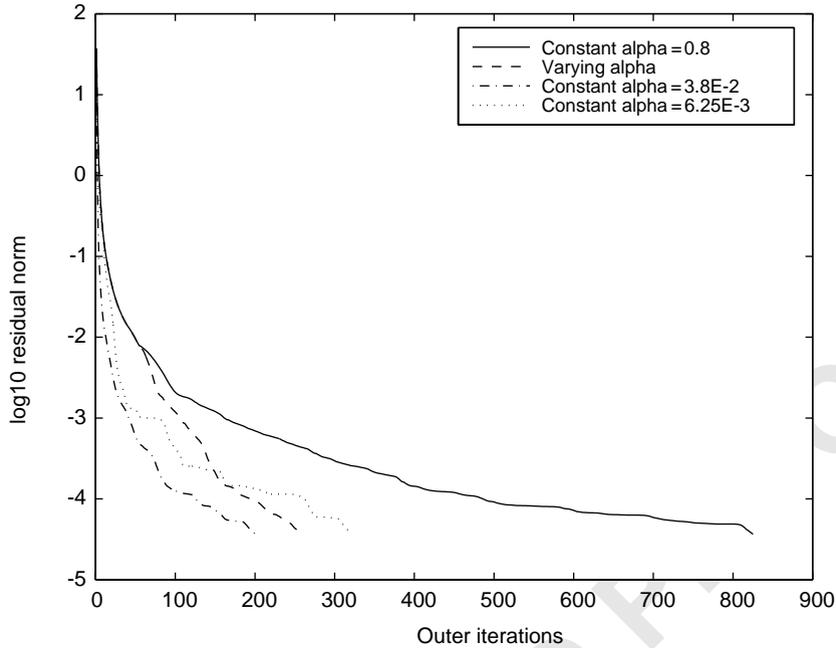


Fig. 8. Choosing α for the problem ELTCOQUE without rational approximation.

1 incomplete LU factors. Relatively inexpensive modifications could potentially be obtained by means
 2 of sparse approximate inverse techniques as mentioned in [5]. Devising an effective procedure for
 3 updating LU factors is beyond the scope of this paper. In the experiments, we re-factor the shifted
 4 matrix A each time a new shift value is taken. Since this procedure is expensive, it is performed
 5 only at a GMRES restart. The dashed line in Fig. 8 indicates that the iterative convergence is much
 6 faster for varying α dynamically than for some constant *large* α (solid line). In this case, it is faster
 7 not only in terms of iteration numbers as shown in Fig. 8 but also in terms of the execution times,
 8 which include refactoring for varying α . The solution times are 222.22 and 303.54 s, respectively.
 9 Monitoring \mathcal{E}_α allows an early detection of a possible preconditioner instability for some small
 10 α indicating that it should not be decreased further. We can also use \mathcal{E}_α to find a better constant
 11 shift and re-construct a preconditioner with this shift. The dash-dotted and dotted curves show the
 12 convergence histories for the two constant shift values ($3.8 \cdot 10^{-2}$ and $6.25 \cdot 10^{-3}$, respectively),
 13 which are obtained at the points immediately preceding a sharp increase in the estimate of \mathcal{E}_α for
 14 two different strategies of decreasing α , fast and gradual, respectively. In particular, a fast decrease,
 15 which consists of halving α at each restart, quickly reaches a shift value corresponding to an unstable
 16 preconditioner. However, the preconditioner re-constructed with the constant shift value $6.25 \cdot 10^{-3}$
 17 obtained from this fast strategy leads to a slower convergence than the one with the shift $3.8 \cdot 10^{-2}$
 18 obtained from a more gradual α decrease. In fact, the constant shift $3.8 \cdot 10^{-2}$ seems to give the best
 19 convergence among the four shift choices. However, finding this optimal shift value is expensive.
 20 The effect of an unstable preconditioning is especially pronounced when the shift value continues
 21 to be halved (dash-dotted curve in Fig. 9), where the residual norm *increases* at about iteration 500.
 Fig. 9 presents the convergence curves of the experiments in which an iterative method attempts

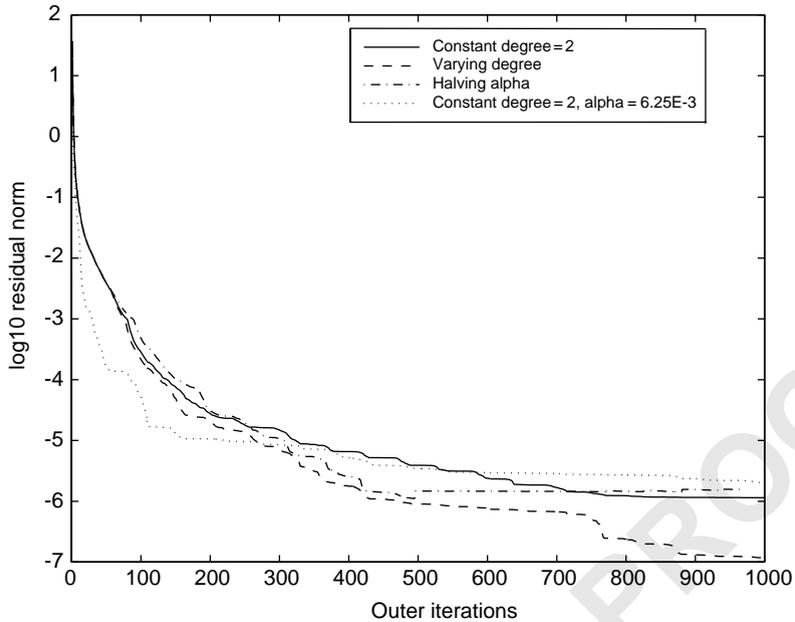


Fig. 9. Achieving high accuracy for the problem ELTCOQUE with rational approximation.

1 to achieve maximum accuracy in 1000 iterations given four ways to choose (α, d) in the rational
 2 approximation (Algorithm 2.2) of the preconditioning. Both the solid and dashed lines are for the
 3 case when alpha is decreasing slowly. The curve for a fixed degree of approximation is represented by
 4 the solid line. The case where the degree is increased by one at each restart, with initial degree 2, is
 5 represented by the dashed line. The dotted line corresponds to the constant smallest shift ($6.25 \cdot 10^{-3}$)
 6 as given in Fig. 8. Note that keeping α constant and small enough accelerates convergence in the
 7 first iterations, but the ultimate residual norm reduction may be much less than when the degree is
 8 varied and a large shift value is taken (dashed curve). Thus varying the degree as well as the shift α
 9 can be beneficial in achieving high accuracy in spite of an increase in the cost of the preconditioning
 10 operation.

11 5. Conclusion

12 We have shown a strategy for building an effective preconditioner for dealing with highly ill-
 13 conditioned matrices. The main difficulty with such matrices is that the standard ILU preconditioners
 14 tend to produce an ILU factorization that is often unstable. Instability can sometimes be avoided
 15 by using a very high level of fill-in to obtain an LU factorization that is very close to that of
 16 A . This approach may not be feasible because of its high memory and computational cost. The
 17 alternative proposed in this paper, is to shift the matrix before computing its ILU factorization, and
 18 then to use a rational expansion to increase the accuracy by extrapolating it to approximate A^{-1} .
 19 We have explained why changing the shift or the degree during iteration helps refocus the iterative
 process in reducing residual components on different parts of the spectrum and can be quite helpful

1 in improving convergence. Numerical experiments support this hypothesis. They also show that the
2 method can succeed in solving rather difficult problems without requiring an excessive amount of
3 memory.

References

- 5 [1] J.L. Batoz, G. Dhatt, Modélisation des structures par éléments finis, Vol. 3: coques, Hermès, Paris, 1992.
6 [2] P.N. Brown, Y. Saad, Hybrid Krylov methods for nonlinear systems of equations, *SIAM J. Sci. Statist. Comput.* 11
7 (1990) 450–481.
8 [3] A. Chapman, Y. Saad, Deflated and augmented Krylov subspace techniques, *Numer. Linear Algebra Appl.* 4 (1997)
9 43–66.
10 [4] E. Chow, Y. Saad, Experimental study of ILU preconditioners for indefinite matrices, *J. Comput. Appl. Math.* 86
11 (1997) 387–414.
12 [5] E. Chow, Y. Saad, Approximate inverse preconditioners via sparse-sparse iterations, *SIAM J. Sci. Comput.* 19 (1998)
13 995–1023.
14 [6] H.C. Elman, A stability analysis of incomplete LU factorizations, *Math. Comput.* 47 (1986) 191–217.
15 [7] R. Glowinski, H.B. Keller, L. Reinhart, Continuation-conjugate gradient methods for the least squares solution of
16 nonlinear boundary value problems, *SIAM J. Sci. Statist. Comput.* 6 (1985) 793–832.
17 [8] T.A. Manteuffel, An incomplete factorization technique for positive definite linear systems, *Math. Comput.* 34 (1980)
18 473–497.
19 [9] R.B. Morgan, A restarted GMRES method augmented with eigenvectors, *SIAM J. Matrix Anal. Appl.* 16 (1995)
20 1154–1171.
21 [10] B.N. Parlett, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, NJ, 1980.
22 [11] A. Ruhe, Rational Krylov sequence methods for eigenvalue computations, *Linear Algebra Appl.* 58 (1984) 391–405.
23 [12] Y. Saad, A flexible inner-outer preconditioned GMRES algorithm, *SIAM J. Sci. Statist. Comput.* 14 (1993)
24 461–469.
25 [13] Y. Saad, Theoretical error bounds and general analysis of a few Lanczos-type algorithms, in: J.D. Brown, M.T.
26 Chu, D.C. Ellison, R.J. Plemmons (Eds.), *Proceedings of the Cornelius Lanczos International Centenary Conference*,
27 SIAM, Philadelphia, PA, 1994, pp. 123–134.
28 [14] Y. Saad, *Iterative Methods for Sparse Linear Systems*, PWS Publishing, New York, 1996.
29 [15] Y. Saad, B. Suchoamel, ARMS: an algebraic recursive multilevel solver for general sparse linear systems, Technical
30 Report umsi-99-107, Minnesota Supercomputer Institute, University of Minnesota, Minneapolis, MN, 1999.
31 [16] M. Sosonkina, J.T. Melson, Y. Saad, L.T. Watson, Preconditioning strategies for linear systems in tire design, *Numer.*
Linear Algebra Appl. 7 (2000) 743–757.